

On the Diffusion of Test Smells in Automatically Generated Test Code: An Empirical Study

Fabio Palomba*, Dario Di Nucci*, Annibale Panichella⁺,
Rocco Oliveto[°], Andrea De Lucia*

A stylized, abstract graphic of an orange flame or smoke, swirling upwards and to the left, occupies the left side of the slide.

SBST 2016
May, 16th 2016
Austin, TX, USA

* University of Salerno, Italy

⁺Delft University of Technology, The Netherlands

[°]University of Molise, Italy



Delft
University of
Technology



Automatic Generation of Test Code



Automatic Generation of Test Code

Effectiveness of Whole Suite Test Case Generation

[Arcuri and Fraser - SSBSE'14]

On The Effectiveness of Whole Test Suite Generation

Effectiveness of Generated Test Cases on Effectiveness

[Shamshiri et al. - ASE'15]

2015 30th IEEE/ACM International Conference on Automated Software Engineering

Impact of Test Case Summaries on Bug Fixing Performance

[Panichella et al. - ICSE'16]

The Impact of Test Case Summaries on Bug Fixing Performance: An Empirical Investigation

Automatic Generation of Test Code

Usability of Automatic Generation Tools in Practice [Rojas et al. - ISSTA'15]

Automated Unit Test Generation during Software Development: A Controlled Experiment and Think-Aloud Observations

José Miguel Rojas¹ Gordon Fraser¹ Andrea Arcuri²

¹Department of Computer Science, The University of Sheffield, United Kingdom
²Scientia, Norway and SNT Centre, University of Luxembourg, Luxembourg

ABSTRACT

Automated unit test generation tools can produce tests that are superior to manually written ones in terms of code coverage, but are these tests helpful to developers while they are writing code? A developer would first need to know when and how to apply such a tool, and would then need to understand the resulting tests in order to provide test oracles and to diagnose and fix any faults that the tests reveal. Considering all this, does automatically generating unit tests provide any benefit over simply writing unit tests manually?

We empirically investigated the effects of using an automated unit test generation tool (EVOSETTE) during development. A controlled experiment with 41 students shows that using EVOSETTE leads to an average branch coverage increase of +13% and 36% less time is spent on testing compared to writing unit tests manually. However, there is no clear effect on the quality of the implementations, as it depends on how the test generation tool and the generated tests are used. In-depth analysis, using five think-aloud observations with professional programmers, confirms the necessity to increase the *usability* of automated unit test generation tools, to *integrate* them better during software development, and to *educate* software developers on how to best use those tools.

Categories and Subject Descriptors: D.2.5 [Software Engineering]: Testing and Debugging – Testing Tools;

Keywords: Automated unit test generation, controlled experiment, think-aloud observations, unit testing

1. INTRODUCTION

Modern automated test generation tools and techniques can efficiently generate program inputs that lead to execution of almost any desired program location. In the context of automated unit test generation for object oriented software there are tools that exercise code contracts [1, 2] or parameterised unit tests [3], try to exhibit undecidable exceptions [1, 4], or simply aim to achieve high coverage [5–9]. There are even commercial tools like Agitar One [10] and Parasoft Test [11] that generate unit tests that capture the current behaviour for automated regression testing. In previous work [12], we have shown that *testers* can achieve higher code coverage when using automatically generated tests than when testing manually. However, most previous investigations considered unit test generation independently of the *developers*: Does the use of an automated unit test generation tool support software developers in *writing code and unit tests*? And if not, how do automated unit test generation tools need to be improved in order to become useful?

In order to provide a better understanding of the effects automatic unit test generation has on software developers, we empirically studied a scenario where developers are given the task of implementing a Java class and accompanying JUnit test suite. We studied this scenario using two empirical methods: First, we performed a controlled empirical study using 41 human subjects, who were asked to complete two coding and testing tasks: one manually, and one assisted by the EVOSETTE unit test generation tool [13]. Second, we conducted five think-aloud observations with professional programmers to understand in detail how developers interact with the testing tool during implementation and testing. Our experiments yielded the following key results:

Effectiveness: Generated unit tests tend to have higher code coverage, and they can support developers in improving the coverage of their manually written tests (on average +7% instruction coverage, +13% branch coverage, +5% mutation score), confirming previous results [12]. However, the influence of automated unit test generation on the quality of the implemented software depends on how a tool and its tests are used.

Efficiency: When using automated unit test generation, developers spend less time on unit testing (36% less time in our experiments). Although participants commented positively on this time saving, it comes at the price of higher uncertainty about the correctness of the implementation, and we observe that spending more time with the generated tests leads to better implementations.

Usability: Readability of automatically generated unit tests is a key aspect that needs to be optimised (63% of participants commented that readability is the most difficult aspect of generated tests). If generated unit tests do not represent understandable, realistic scenarios, they may even have detrimental effects on software.

Integration: Test generation tools typically optimise for code coverage, but this is not what developers do. We observed three different approaches to unit testing, each posing different requirements on how testing tools should interact with developers and existing tests.

Education: Simply providing a test generation tool to a developer will not necessarily lead to better software—developers need education and experience on when and how to use such tools, and we need to establish best practices. For example, we found a moderate correlation between the frequency of use of our test generation tool and the number of implementation errors committed by developers.

The main contributions of this paper are the controlled empirical study (Section 3) and the think-aloud observations (Section 4).

Copyright is held by the owner(s). Publication rights licensed to ACM.
ISSTA'15, July 13–17, 2015, Baltimore, MD, USA
ACM, 978-1-4503-3620-4/15/07
<http://dx.doi.org/10.1145/2771783.2771801>

978-1-4503-3620-4/15/07 \$15.00 © 2015 ACM, Inc. This is the author's version of the work. It is posted here by permission of ACM for your personal use. Not for redistribution. The definitive version was published in ISSTA'15, July 13–17, 2015, Baltimore, MD, USA. ACM, 978-1-4503-3620-4/15/07
Copyright is held by the owner(s). Publication rights licensed to ACM.
ISSTA'15, July 13–17, 2015, Baltimore, MD, USA
ACM, 978-1-4503-3620-4/15/07
<http://dx.doi.org/10.1145/2771783.2771801>

Automatic Generation of Test Code



What about the characteristics of test code produced by such tools?

Test Smells in Test Code

Refactoring Test Code

Arie van Deursen Leon Moonen
CWI
The Netherlands
<http://www.cwi.nl/~{arie,leon}/>
{arie,leon}@cwi.nl

Alex van den Bergh Gerard Kok
Software Improvement Group
The Netherlands
<http://www.software-improvers.com/>
{alex,gerard}@software-improvers.com

ABSTRACT
Two key aspects of extreme programming (XP) are unit testing and merciless refactoring. Given the fact that the ideal test code / production code ratio approaches 1:1, it is not surprising that unit tests are being refactored. We found that refactoring test code is different from refactoring production code in two ways: (1) there is a distinct set of bad smells involved, and (2) improving test code involves additional test-specific refactorings. To share our experiences with other XP practitioners, we describe a set of bad smells that indicate trouble in test code, and a collection of test refactorings to remove these smells.

Keywords
Refactoring, unit testing, extreme programming.

1 INTRODUCTION

"If there is a technique at the heart of *extreme programming* (XP), it is unit testing" [1]. As part of their programming activity, XP developers write and maintain (white box) unit tests continually. These tests are automated, written in the same programming language as the production code, considered an explicit part of the code, and put under revision control.

The XP process encourages writing a test class for every class in the system. Methods in these test classes are used to verify complicated functionality and unusual circumstances. Moreover, they are used to document code by explicitly indicating what the expected results of a method should be for typical cases. Last but not least, tests are added upon receiving a bug report to check for the bug and to check the bug fix [2]. A typical test for a particular method includes: (1) code to set up the fixture (the data used for testing), (2) the call of the method, (3) a comparison of the actual results with the expected values, and (4) code to tear down the fixture. Writing tests is usually supported by frameworks such as *JUnit* [3].

The test code / production code ratio may vary from project to project, but is ideally considered to approach a ratio of 1:1. In our project we currently have a 2:3 ratio, although others have reported a lower ratio¹. One of the corner stones of XP is that having many tests available helps the developers to overcome their fear for change: the tests will provide immediate feedback if the system gets broken at a critical place. The downside of having many tests, however, is that changes in functionality will typically involve changes in the test code as well. The more test code we get, the more important it becomes that this test code is as easily modifiable as the production code.

The key XP practice to keep code flexible is "refactor mercilessly": transforming the code in order to bring it in the simplest possible state. To support this, a catalog of "code smells" and a wide range of refactorings is available, varying from simple modifications up to ways to introduce design patterns systematically in existing code [5].

When trying to apply refactorings to the test code of our project we discovered that refactoring test code is different from refactoring production code. Test code has a distinct set of smells, dealing with the ways in which test cases are organized, how they are implemented, and how they interact with each other. Moreover, improving test code involves a mixture of refactorings from [5] specialized to test code improvements, as well as a set of additional refactorings, involving the modification of test classes, ways of grouping test cases, and so on.

The goal of this paper is to share our experience in improving our test code with other XP practitioners. To that end, we describe a set of *test smells* indicating trouble in test code, and a collection of *test refactorings* explaining how to overcome some of these problems through a simple program modification.

This paper assumes some familiarity with the xUnit framework [3] and refactorings as described by Fowler [5]. We will refer to refactorings described in this book using *Name*

¹ This project started a year ago and involves the development of a product called DocGen [4]. Development is done by a small team of five people using XP techniques. Code is written in Java and we use the JUnit framework for unit testing.

"Test Smells represent a set of a poor design solutions to write tests "

[Van Deursen et al. - XP 2001]



test smells related to the way
developers write test fixtures
and test cases

Test Smells in Test Code

```
public void test12 () throws Throwable {  
    JSTerm jSTerm0 = new JSTerm();  
    jSTerm0.makeVariable () ;  
    jSTerm0.add((Object) "");  
    jSTerm0.matches(jSTerm0);  
    assertEquals (false, jSTerm0.isGround ());  
    assertEquals(true, jSTerm0.isVariable());  
}
```

Test Smells in Test Code

```
public void test12 () throws Throwable {  
    JSTerm jSTerm0 = new JSTerm();  
    jSTerm0.makeVariable () ;  
    jSTerm0.add((Object) "");  
    jSTerm0.matches(jSTerm0);  
    assertEquals (false, jSTerm0.isGround ());  
    assertEquals(true, jSTerm0.isVariable());  
}
```

The test method **checks** the production method `isGround()`

Test Smells in Test Code

```
public void test12 () throws Throwable {  
    JSTerm jSTerm0 = new JSTerm();  
    jSTerm0.makeVariable () ;  
    jSTerm0.add((Object) "");  
    jSTerm0.matches(jSTerm0);  
    assertEquals (false, jSTerm0.isGround ());  
    assertEquals(true, jSTerm0.isVariable());  
}
```

But also the production method isVariable()

Test Smells in Test Code

```
public void test12 () throws Throwable {  
    JSTerm jSTerm0 = new JSTerm();  
    jSTerm0.makeVariable () ;  
    jSTerm0.add((Object) "");  
    jSTerm0.matches(jSTerm0);  
    assertEquals (false, jSTerm0.isGround ());  
    assertEquals(true, jSTerm0.isVariable());  
}
```

This is an Eager Test, namely a test which checks more than one method of the class to be tested, making difficult the comprehension of the actual test target.

Test Smells in Test Code

A test case is affected by a Resource Optimism when it makes assumptions about the state or the existence of external resources, providing a non-deterministic result that depend on the state of the resources.

An Assertion Roulette comes from having a number of assertions in a test method that have no explanation. If an assertion fails, the identification of the assert that failed can be difficult.

Who cares about Test Smells?
Test Cases can be re-generated!



Who cares about Test Smells?
Test Cases can be re-generated!

True



Who cares about Test Smells?
Test Cases can be re-generated!

True

BUT



Who cares about Test Smells?

Test Cases can be re-generated!

True

BUT

Developers modify and remove test code

Developers add tests when automatic tools leave uncovered branches

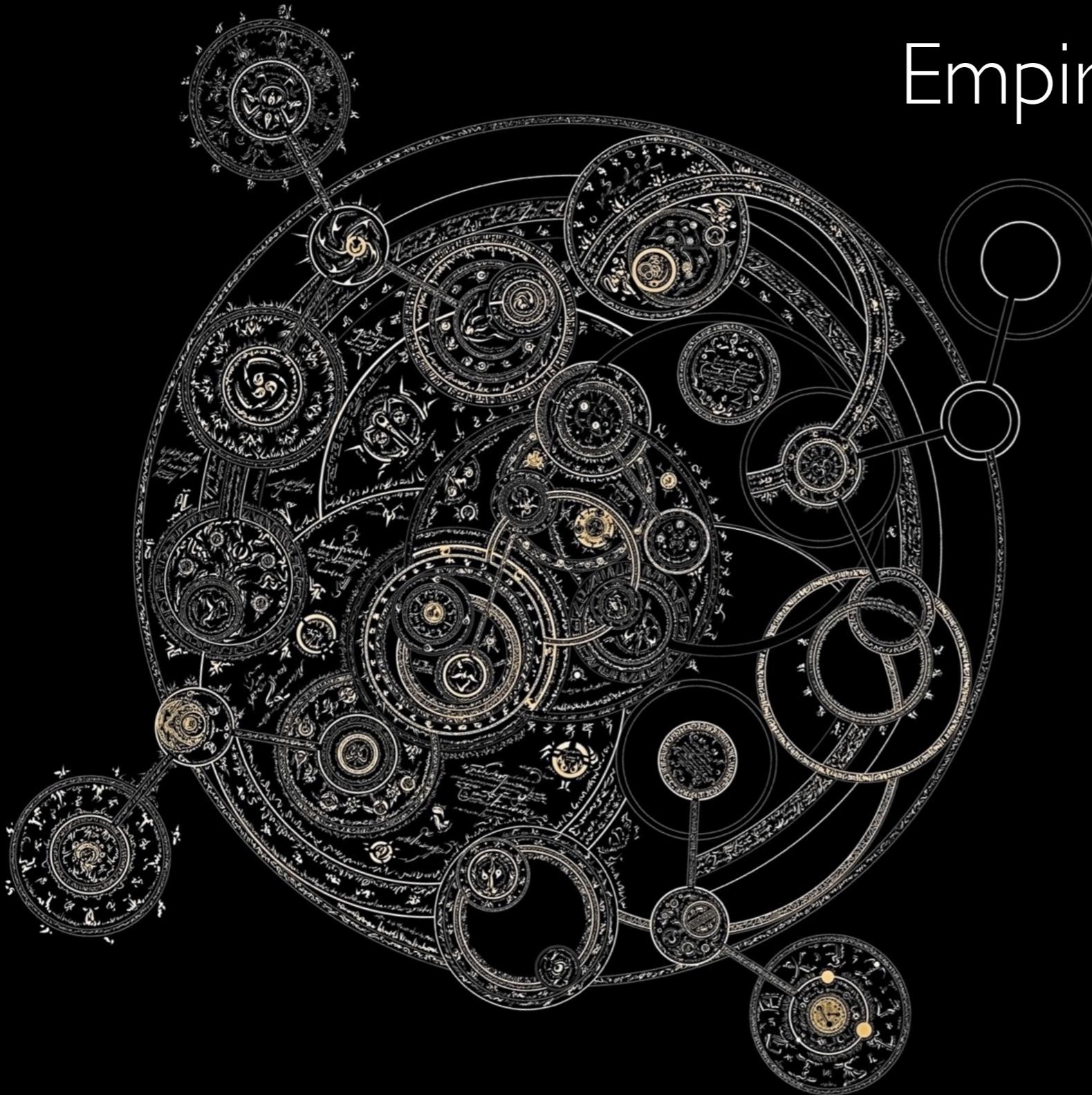
Developers combine generated with manually written tests

Usability of Automatic Generation Tools in Practice

[Rojas et al. - ISSTA'15]

On the Diffusion of Test Smells in Automatically Generated Test Code: An Empirical Study

Empirical Study Design



Empirical Study Design

8

test smell types

“Refactoring Test Code”

[Van Deursen et al. - XP 2001]

Empirical Study Design

8

test smell types

“Refactoring Test Code”
[Van Deursen et al. - XP 2001]

110

software projects

“A Large Scale Evaluation of Automated Unit
Test Generation using Evosuite”
[Fraser and Arcuri - TOSEM 2014]

Empirical Study Design

8

test smell types

“Refactoring Test Code”
[Van Deursen et al. - XP 2001]

110

software projects

“A Large Scale Evaluation of Automated Unit
Test Generation using Evosuite”
[Fraser and Arcuri - TOSEM 2014]

EVOSUITE

Empirical Study Design

16,603

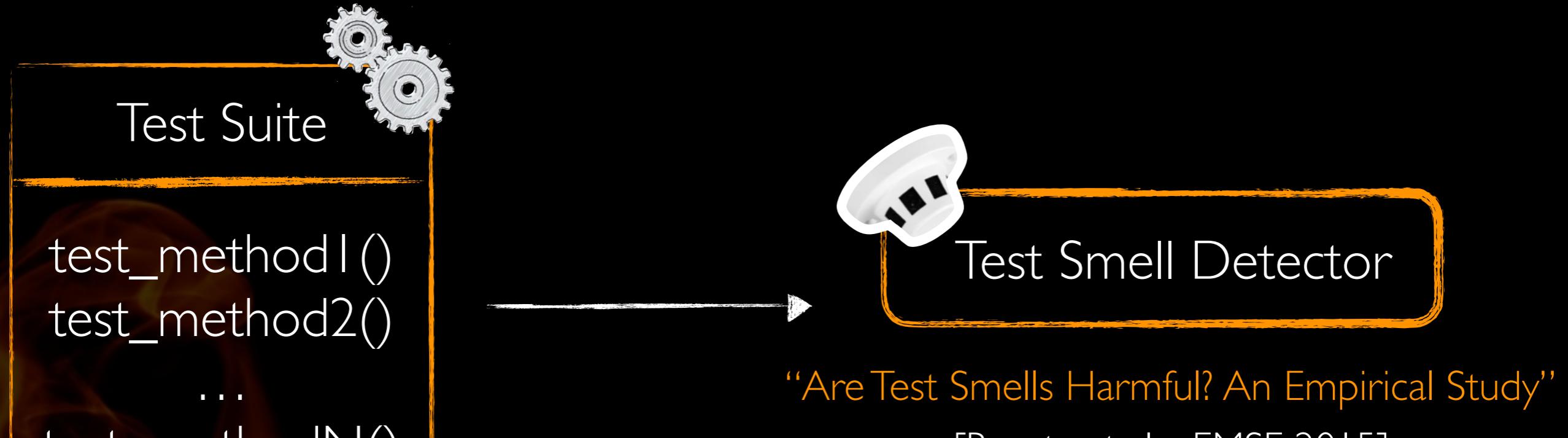
JUnit classes

“A Large Scale Evaluation of Automated Unit
Test Generation using Evosuite”

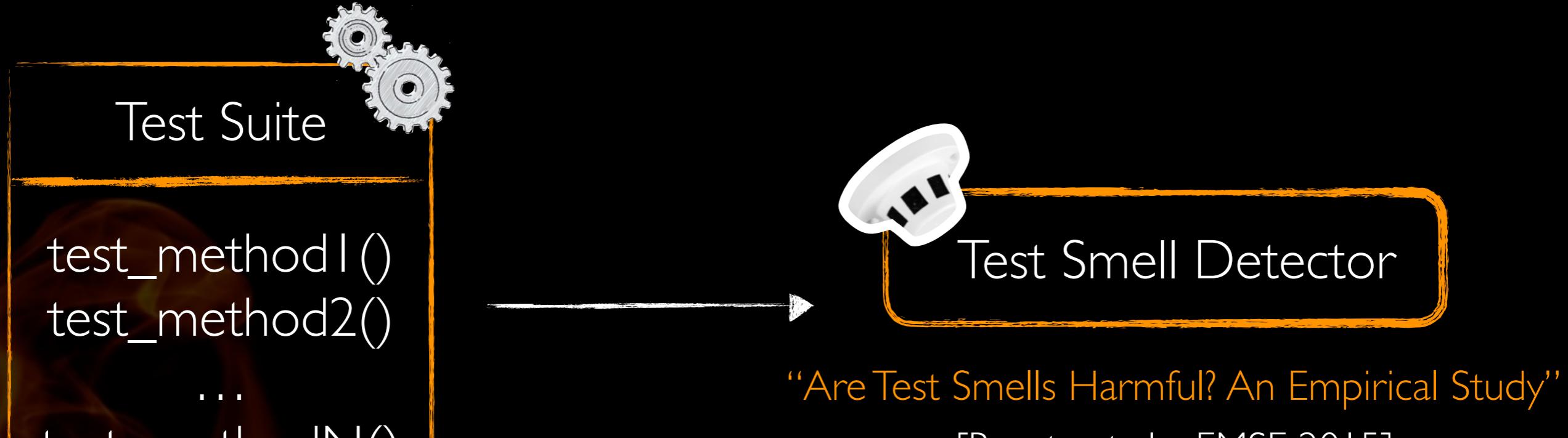
[Fraser and Arcuri - TOSEM 2014]



Data Extraction



Data Extraction



75% 100%
precision recall

Sample size: 378 JUnit classes

Research Questions

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

Research Questions

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Research Questions

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

RQ3: Which Test Smells Co-Occur Together?

Research Questions

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

RQ3: Which Test Smells Co-Occur Together?

RQ4: Is There a Relationship Between the Presence of Test Smells and the Project Characteristics?

On the Diffusion of Test Smells in Automatically Generated Test Code: An Empirical Study

Analysis of the Results



Results of the Study

RQ1: To What Extent Test Smells are Spread in
Automatically Generated Test Classes?

13,791

smelly JUnit classes

Results of the Study

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

83%

of the JUnit classes analyzed

Results of the Study

RQ1: To What Extent Test Smells are Spread in Automatically Generated Test Classes?

RQ1

Test Smells are highly diffused in the automatically generated test suites

Results of the Study

RQ2: Which Test Smells Occur More Frequently
in Automatically Generated Test Classes?

Assertion Roulette

54%

Results of the Study

RQ2: Which Test Smells Occur More Frequently
in Automatically Generated Test Classes?

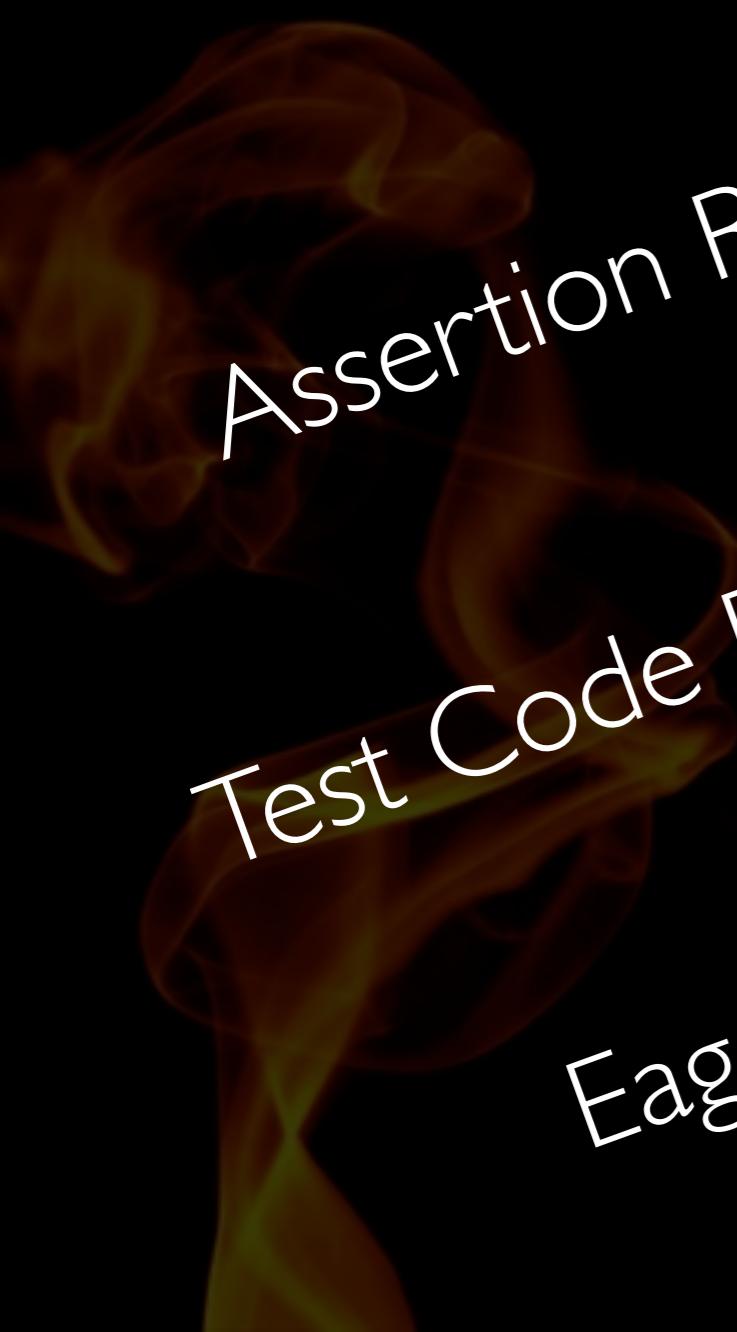
Assertion Roulette
Test Code Duplication

54%

33%

Results of the Study

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

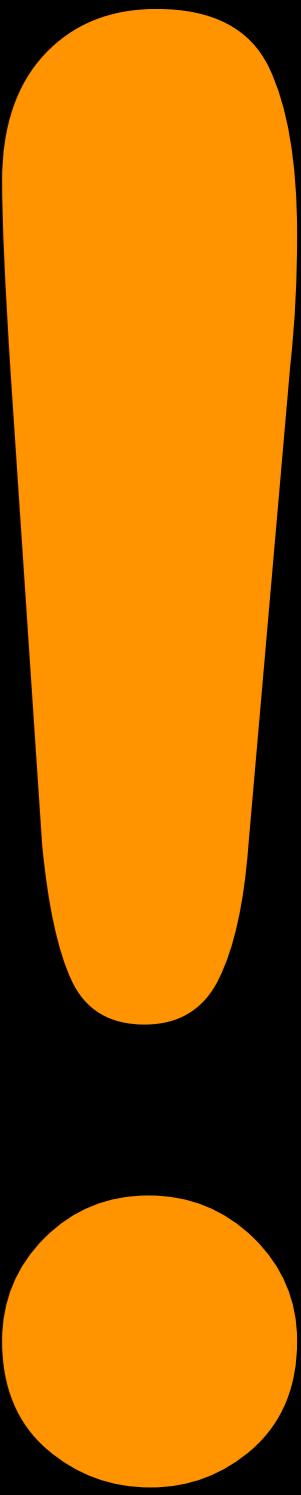


Assertion Roulette
Test Code Duplication
Eager Test

54%

33%

29%



Results of the Study

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Assertion Roulette

```
public void test8 () throws Throwable {  
    Document document0 = new Document();  
    assertNotNull(document0);  
  
    document0.procText.add((Character) "s");  
  
    String string0 = document0.stringify();  
    assertEquals ("s", document0.stringify());  
    assertNotNull(string0);  
    assertEquals("s", string0);  
}
```

Results of the Study

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Assertion Roulette

What is the behavior under test?

Are the generated assertions valid?

Results of the Study

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Assertion Roulette

These problems have a huge impact on developers' ability to find faults

The Impact of Test Case Summaries on Bug Fixing Performance
An Empirical Investigation
[Panichella et al. - ICSE'16]

Results of the Study

RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Test Code Duplication

```
public void test8 () throws Throwable {  
    GenericProperties generic0 = new GenericProperties();  
    boolean boolean0 = generic0.isValidClassname();  
    ...  
}
```

```
public void test9 () throws Throwable {  
    GenericProperties generic0 = new GenericProperties();  
    boolean boolean0 = generic0.isValidClassname();  
    ...  
}
```

Results of the Study

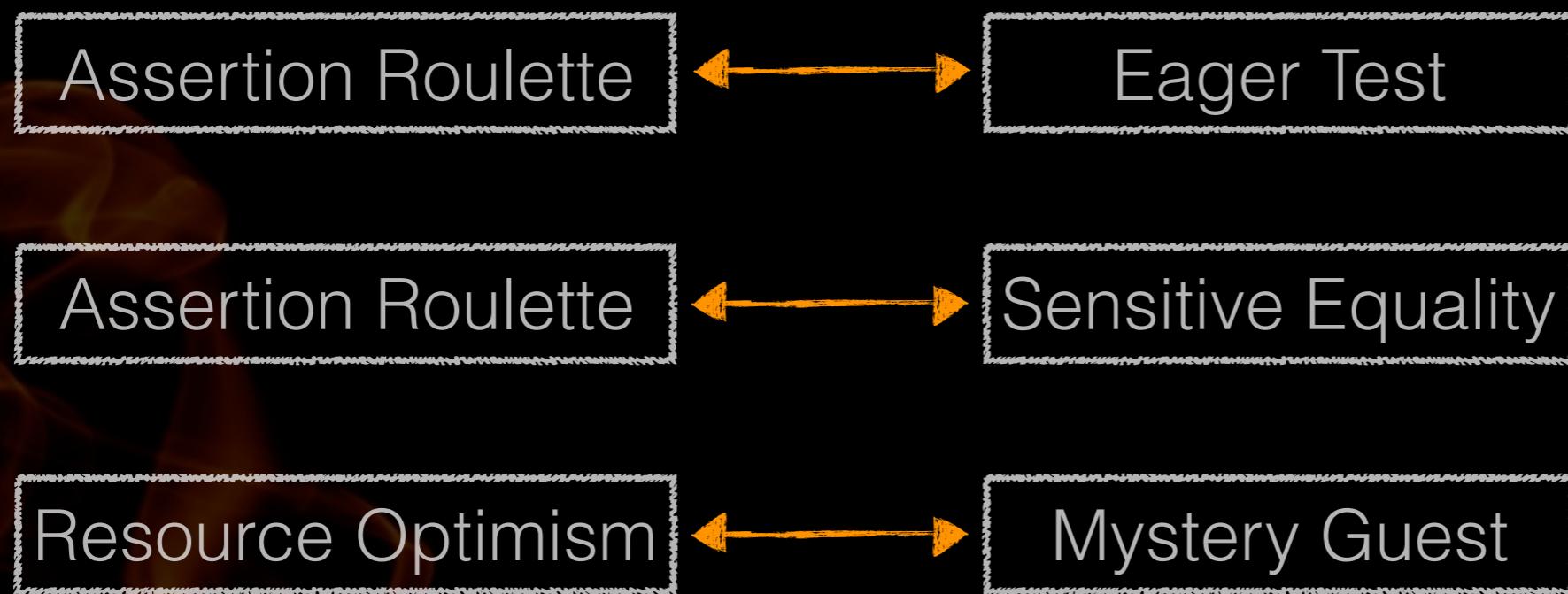
RQ2: Which Test Smells Occur More Frequently in Automatically Generated Test Classes?

Test Code Duplication

This problem can be avoided by generating test fixtures!

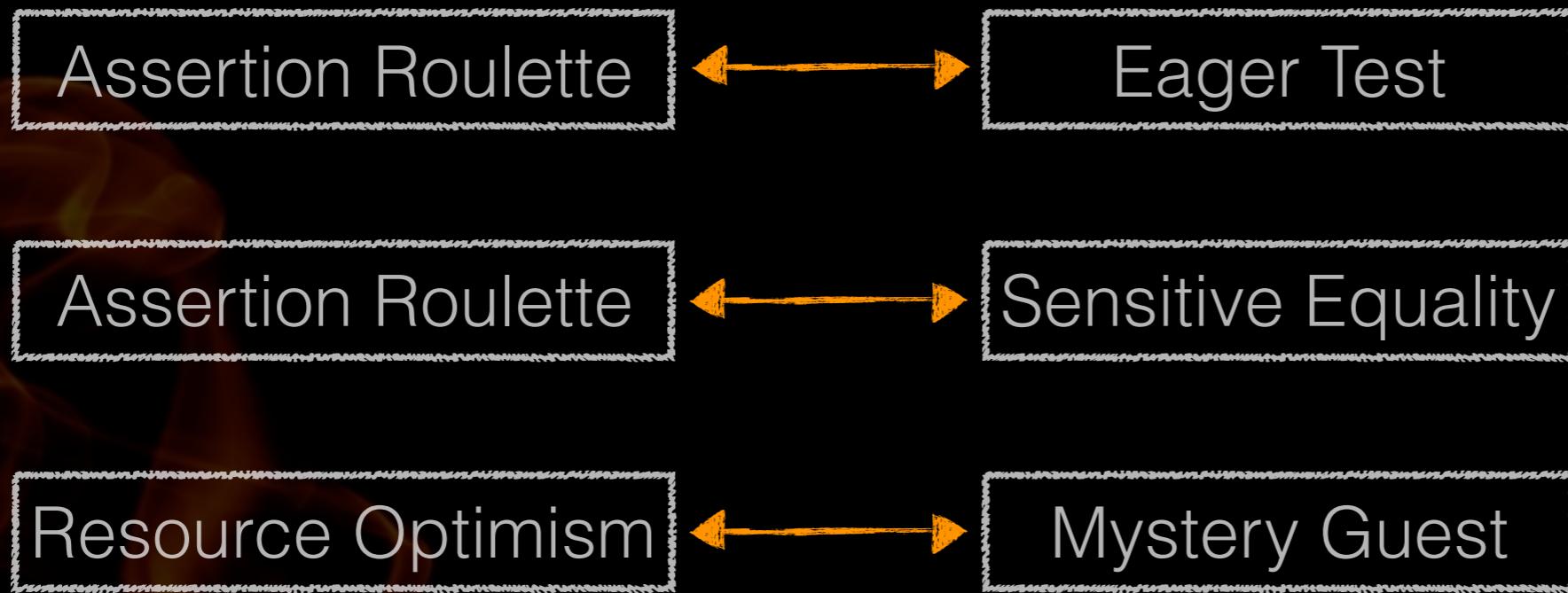
Results of the Study

RQ3: Which Test Smells Co-Occur Together?



Results of the Study

RQ3: Which Test Smells Co-Occur Together?



Automatic tools have as main goal that of maximize coverage, without considering test code quality

Results of the Study

RQ4: Is There a Relationship Between the Presence of Test Smells and the Project Characteristics?

Yes! The higher the LOC to be tested, the higher the probability to produce a smelly test case!

Results of the Study

RQ4: Is There a Relationship Between the Presence of Test Smells and the Project Characteristics?

Yes! The higher the LOC to be tested, the higher the probability to produce a smelly test case!

The higher the LOC of the JUnit class, the higher the probability to produce a smelly test case!

Summarizing

Lesson I: Current implementations of search-based algorithms for automatic test case generation do not consider code quality, increasing the probability to introduce smells!

Summarizing

Lesson I: Current implementations of search-based algorithms for automatic test case generation do not consider code quality, increasing the probability to introduce smells!

NB: Considering test code quality is important not only to avoid the introduction of smells, but also because the coverage can be increased!

Automatic Test Case Generation: What If Test Code Quality Matters?
[F. Palomba, A. Panichella, A. Zaidman, R. Oliveto, A. De Lucia - ISSTA'16]

Summarizing

Lesson 2: Automatic test case generation tools do not produce text fixtures during their computation, and this implies the introduction of several code clones in the resulting JUnit classes.

Future research should spend effort in the automatic generation of test fixtures!

From now on...

Challenge I: Evaluating Test Smells in Test Cases
automatically generated by other tools

From now on...

Challenge 1: Evaluating Test Smells in Test Cases
automatically generated by other tools

Challenge 2: Defining new algorithms able to solve
the design problems analyzed (e.g., test fixtures).

On the Diffusion of Test Smells in Automatically Generated Test Code: An Empirical Study

Fabio Palomba*, Dario Di Nucci*, Annibale Panichella⁺,
Rocco Oliveto[°], Andrea De Lucia*

A stylized, abstract graphic of an orange flame or smoke, swirling upwards and to the left, occupies the left side of the slide.

SBST 2016
May, 16th 2016
Austin, TX, USA

* University of Salerno, Italy

⁺Delft University of Technology, The Netherlands

[°]University of Molise, Italy